

Statistics for Social and Behavioral Sciences

Michael J. Kolen
Robert L. Brennan

Test Equating, Scaling, and Linking

Methods and Practices

Third Edition

 Springer

Statistics for Social and Behavioral Sciences

Series editors

S. E. Fienberg

W. J. van der Linden

For further volumes:

<http://www.springer.com/series/3463>

Michael J. Kolen · Robert L. Brennan

Test Equating, Scaling, and Linking

Methods and Practices

Third Edition

 Springer

Michael J. Kolen
Iowa Testing Programs
University of Iowa
Iowa City, IA
USA

Robert L. Brennan
CASMA
University of Iowa
Iowa City, IA
USA

ISBN 978-1-4939-0316-0 ISBN 978-1-4939-0317-7 (eBook)
DOI 10.1007/978-1-4939-0317-7
Springer New York Heidelberg Dordrecht London

Library of Congress Control Number: 2013956631

© Springer Science+Business Media New York 2014

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law. The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

To Amy, Raychel, and Daniel

—M. J. K.

To Cicely and Sean

—R. L. B.

Preface

Prior to 1980, the subject of equating was ignored by most people in the measurement community except for psychometricians, who had responsibility for equating. Beginning in the early 1980s, the importance of equating was recognized by a broader spectrum of people associated with testing. This increased attention to equating is attributable to at least three developments. First, there continues to be an increase in the number and variety of testing programs that use multiple forms of tests, and the testing professionals responsible for such programs have recognized that scores on multiple forms should be equated. Second, test developers and publishers often have referenced the role of equating in arriving at reported scores to address a number of issues raised by testing critics. Third, the accountability movement in education and issues of fairness in testing have become much more visible. These developments have given equating an increased emphasis among measurement professionals and test users.

In addition to statistical procedures, successful equating involves many aspects of testing, including procedures to develop tests, to administer and score tests, and to interpret scores earned on tests. Of course, psychometricians who conduct equating need to become knowledgeable about all aspects of equating. The prominence of equating, along with its interdependence with so many aspects of the testing process, also suggests that test developers and all other testing professionals should be familiar with the concepts, statistical procedures, and practical issues associated with equating.

Before we published the first edition in 1995, the need for a book on equating became evident to us from our experiences in equating hundreds of test forms in many testing programs, in training psychometricians to conduct equating, in conducting seminars and courses on equating, and in publishing on equating and other areas of psychometrics. Our experience suggested that relatively few measurement professionals had sufficient knowledge to conduct equating. Also, many did not fully appreciate the practical consequences of various changes in testing procedures on equating, such as the consequences of many test-legislation initiatives, the use of constructed-response items in assessments, and the introduction of computer-based test administration. Consequently, we believed that measurement professionals needed to be educated in equating methods and practices; the 1995 book was intended to help fulfill this need. Although several general published references on equating existed at the time (e.g., Angoff 1971;

Harris and Crouse 1993; Holland and Rubin 1982; Petersen et al. 1989), none of them provided the broad, integrated, in-depth, and up-to-date coverage of the first edition of this book.

After the publication of the first edition in 1995, a large body of new research was published. Much of this work was in technical areas that include smoothing in equipercentile equating, estimation of standard errors of equating, and the use of polytomous item response theory methods in equating. In addition, the use of constructed-response items and computer-based tests became more prominent. These applications create complexities for equating beyond what is typically encountered with paper-and-pencil multiple-choice tests. Thus, updating the material in the first edition was one of the reasons for publishing a second edition.

The first edition briefly considered score scales and test linking. The second edition devoted whole chapters to each of these topics. The development of score scales is an important component of the scaling and equating process. Linking of tests has been of much recent interest, due to various investigations of how to link tests from different test publishers or constructed for different purposes (e.g., Feuer et al. 1999). Because both scaling and linking are closely related to test equating, it seemed natural to extend coverage along these lines.

Following the publication of the second edition in 2004, a considerable amount of research was conducted on equating, scaling, and linking. In addition to a substantial number of journal articles, Dorans, Pommerich, and Holland (2007) and von Davier (2011) published edited books on equating, scaling, and linking. In addition, a substantial chapter by Holland and Dorans (2006) provides a conceptual framework for classifying equating and linking methodology that focuses on the properties of scores that are linked and on the requirements of different types of linking. A chapter by Kolen (2006) provides an updated discussion of score scales. The third edition updates all chapters to incorporate this recent literature. Following is a brief overview of the chapters of the third edition.

In [Chap. 1](#), a general introduction is provided, primarily in terms of a conceptual overview. In this chapter, we define equating, describe its relationship to test development, and distinguish equating from scaling and linking. We also present equating designs, properties of equating, and introduce the concept of equating error.

In [Chap. 2](#), using the random groups design, we illustrate traditional equating methods, such as equipercentile and linear methods. We also discuss here many of the key concepts of equating, such as properties of converted scores and the influence of the resulting scale scores on the choice of an equating result.

In [Chap. 3](#), we cover smoothing methods in equipercentile equating. We show that the purpose of smoothing is the reduction of random error in estimating equating relationships in the population. We describe methods based on log-linear models, cubic splines, and strong true score models.

In [Chap. 4](#), we treat linear equating with nonequivalent groups of examinees. We derive statistical methods and stress the need to disconfound examinee-group and test-form differences. Also, we distinguish observed score equating from true score equating.

In [Chap. 5](#), we continue our discussion of equating with nonequivalent groups with a presentation of equipercentile methods.

In [Chap. 6](#), we describe item response theory (IRT) equating methods under various designs. This chapter covers issues that include scaling person and item parameters, IRT true and observed score equating methods, equating using item pools, and equating using polytomous IRT models.

[Chapter 7](#) focuses on standard errors of equating; both bootstrap and analytic procedures are described. We illustrate the use of standard errors to choose sample sizes for equating and to compare the precision in estimating equating relationships for different designs and methods.

In [Chap. 8](#), we describe many practical issues in equating, including the importance of test development procedures, test standardization conditions, and quality control procedures. We stress conditions that are conducive to adequate equating. Also, we discuss comparability issues for mixed-format assessments and computer-based tests.

[Chapter 9](#) is devoted to score scales for tests. We discuss different scaling perspectives. We describe linear and nonlinear transformations that are used to construct score scales, and we consider procedures for enhancing the meaning of scale scores that include incorporating normative, content, and score precision information. We discuss procedures for maintaining score scales and scales for batteries and composites. We conclude with a section on vertical scaling that includes consideration of scaling designs and psychometric methods and a review of research on vertical scaling.

In [Chap. 10](#), we describe linking categorization schemes and criteria and consider equating, vertical scaling, and other related methodologies as a part of these categorization schemes. An extensive example is used to illustrate how the lack of group invariance in concordance relationships can be examined and used as a means for demonstrating some of the limitations of linking methods.

We use a random groups illustrative equating example in [Chaps. 2, 3, and 7](#); a nonequivalent groups illustrative example in [Chaps. 4–6](#); a second random groups illustrative example in [Chaps. 6 and 9](#); and a single-group illustrative example in [Chap. 10](#). We use data from the administration of a test battery in multiple grades for an illustrative example in [Chap. 9](#), and data from the administration of two different tests for an illustrative example in [Chap. 10](#). [Chapters 1–10](#) each have a set of exercises that are intended to reinforce the concepts and procedures in the chapter. The answers to the exercises are in [Appendix A](#). We describe computer programs and how to obtain them in [Appendix B](#).

In addition to updating the review of literature for all of the chapters, the third edition incorporates substantial new material as follows:

- [Chapter 3](#) includes additional procedures to choose models in log-linear pre-smoothing and includes a new brief section on the kernel method of equating.
- [Chapter 4](#) includes a new section on chained linear equating and incorporates chained linear equating in the illustrative example. In addition, it includes a new

discussion of the relationships among linear methods in the common-item nonequivalent groups design.

- **Chapter 5** includes new descriptions of modified frequency estimation equating and chained equipercentile equating, and incorporates these methods in the illustrative example.
- **Chapter 8** includes a new extensive section on equating criteria in research studies. Material on equating mixed-format tests containing multiple-choice and constructed-response items is significantly updated.
- **Chapter 9** includes a new section on unit scores, item scores, and raw scores. A new section on scores for mixed-format tests, including issues in weighting scores for different item types, is added. In addition, a new section on score scales and growth is added.
- **Chapter 10** includes a new summary of the Holland and Dorans (2006) linking framework.

In addition, each chapter contains a reference list, rather than having a single reference list at the end of the volume as in the first two editions.

We anticipate that many readers of this book will be advanced graduate students, entry-level professionals, or persons preparing to conduct equating, scaling, or linking for the first time. Other readers likely will be experienced professionals in measurement and related fields who will want to use this book as a reference. To address these varied audiences, we make frequent use of examples and stress conceptual issues. This book is not a traditional statistics text. Instead, it is meant for instructional use and as a reference for practical use that is intended to address both statistical and applied issues. The most frequently used methodologies are treated, as well as many practical issues. Although we are unable to cover all of the literature on equating, scaling, and linking, we provide many references so that the interested reader may pursue topics of particular interest.

The principal goals of this book are for the reader to understand the principles of equating, scaling, and linking; to be able to conduct equating, scaling, and linking; and to interpret the results in reasonable ways. After studying this book, the reader should be able to

- Understand the purposes of equating, scaling, and linking and the context in which they are conducted.
- Distinguish between equating, scaling, and linking methodologies and procedures.
- Appreciate the importance to equating of test development and quality control procedures.
- Understand the distinctions among equating properties, equating designs, and equating methods.
- Understand fundamental concepts—including designs, methods, errors, and statistical assumptions.
- Compute equating, scaling, and linking functions and choose among methods.
- Interpret results from equating, scaling, and linking analyses.

- Design reasonable and useful equating, scaling, and linking studies.
- Conduct equating, scaling, and linking in realistic testing situations.
- Identify appropriate and inappropriate uses and interpretations of equating, scaling, and linking results.

We cover nearly all of the material in this book in a three semester-hour graduate seminar at The University of Iowa. In our course, we supplement the materials here with general references (Angoff 1971; Holland and Dorans 2006; Holland and Rubin 1982; Petersen et al. 1989) so that the students become familiar with other perspectives and notational schemes.

We have used much of the material in this book in various training sessions, including those at the annual meetings of the National Council on Measurement in Education, the American Educational Research Association, and the American Psychological Association, and in workshops given in Israel, Japan, South Korea, Spain, Taiwan, and The University of Iowa.

We acknowledge the generous contributions that others made to the first edition of this book. We benefitted from interactions with very knowledgeable psychometricians at ACT and elsewhere, and many of the ideas in this book came from conversations and interactions with these people. Specifically, Bradley Hanson reviewed the entire manuscript and made valuable contributions, especially to the statistical presentations. He conducted the bootstrap analyses that are presented in Chapter 7 and, along with Lingjia Zeng, developed much of the computer software used in the examples. Deborah Harris reviewed the entire manuscript, and we thank her especially for her insights on practical issues in equating. Chapters 1 and 8 benefitted considerably from her ideas and counsel. Lingjia Zeng also reviewed the entire manuscript and provided us with many ideas on statistical methodology, particularly in the areas of standard errors and IRT equating. Thanks to Dean Colton for his thorough reading of the entire manuscript, Xiaohong Gao for her review and for working through the exercises, and Ronald Cope and Tianqi Han for reading portions of the manuscript. We are grateful to Nancy Petersen for her in-depth review of a draft of the first edition, her insights, and her encouragement. Bruce Bloxom provided valuable feedback, as did Barbara Plake and her graduate class at the University of Nebraska–Lincoln. We thank an anonymous reviewer, and the reviewer’s graduate student, for providing us with their valuable critique. We are indebted to all who have taken our equating courses and training sessions.

For the second edition, we are grateful to Ye Tong for the many hours she spent on electronic typesetting, for all of the errors she found, and for helping with many of the examples and the exercises. We thank Amy Hendrickson for helping to develop the polytomous IRT examples in Chapter 6, Seonghoon Kim for reviewing the additions to Chapter 6 on polytomous IRT and for developing the computer program POLYST, and Ping Yin for her work on Chapters 4 and 10. We acknowledge the work of Zhongmin Cui and Yueh-Mei Chien on the computer programs, and the work of Noo Ree Huh on checking references. We thank the students in our equating and scaling classes at The University of Iowa who discovered many errors and for helping us clarify some confusing portions of earlier drafts. We are grateful to Neil

Dorans, Samuel Livingston, and Paul Holland for reviewing portions of the new material in the second edition. We express our appreciation to the Iowa Measurement Research Foundation for providing support to the graduate students who worked with us on the second edition. For the third edition, we thank Wei Wang for her many hours spent on electronic typesetting. We also thank many graduate students at The University of Iowa for helping us correct errors that appeared in the second edition. Amy Kolen deserves thanks for her superb editorial advice for all three editions.

Iowa City, IA November, 2013

Michael J. Kolen
Robert L. Brennan

References

- Angoff, W. H. (1971). Scales, norms, and equivalent scores. In R. L. Thorndike (Ed.), *Educational measurement* (2nd ed., pp. 508–600). Washington, DC: American Council on Education.
- Dorans, N. J., Pommerich, M., & Holland, P. (Eds.). (2007). *Linking and aligning scores and scales*. New York: Springer.
- Feuer, M. J., Holland, P. W., Green, B. F., Bertenthal, M. W., & Hemphill, F. C. (Eds.). (1999). *Uncommon measures: Equivalence and linkage among educational tests*. Washington, DC: National Research Council.
- Harris, D. J., & Crouse, J. D. (1993). A study of criteria used in equating. *Applied Measurement in Education*, 6, 195–240.
- Holland, P. W., & Dorans, N. J. (2006). Linking and equating. In R. L. Brennan (Ed.), *Educational measurement* (4th ed., pp. 187–220). Westport, CT: American Council on Education and Praeger.
- Holland, P. W., & Rubin, D. B. (1982). *Test equating*. New York: Academic.
- Kolen, M. J. (2006). Scaling and norming. In R. L. Brennan (Ed.), *Educational measurement* (4th ed., pp. 155–186). Westport, CT: American Council on Education and Praeger.
- Petersen, N. S., Kolen, M. J., & Hoover, H. D. (1989). Scaling, norming, and equating. In R. L. Linn (Ed.), *Educational measurement* (3rd ed., pp. 221–262). New York: Macmillan.
- von Davier, A. A. (Ed.). (2011). *Statistical models for test equating, scaling, and linking*. New York: Springer.

Contents

1	Introduction and Concepts	1
1.1	Equating and Related Concepts	1
1.1.1	Test Forms and Test Specifications	2
1.1.2	Equating	2
1.1.3	Processes That are Related to Equating	3
1.1.4	Equating and Score Scales	4
1.1.5	Equating and the Test Score Decline of the 1960s and 1970s	6
1.2	Equating and Scaling in Practice: A Brief Overview of This Book	7
1.3	Properties of Equating	8
1.3.1	Symmetry Property	9
1.3.2	Same Specifications Property	9
1.3.3	Equity Properties	9
1.3.4	Observed Score Equating Properties	11
1.3.5	Group Invariance Property	12
1.4	Equating Designs	12
1.4.1	Random Groups Design	13
1.4.2	Single Group Design	14
1.4.3	Single Group Design with Counterbalancing	14
1.4.4	ASVAB Problems with a Single Group Design	16
1.4.5	Common-Item Nonequivalent Groups Design	18
1.4.6	NAEP Reading Anomaly: Problems with Common Items	20
1.5	Error in Estimating Equating Relationships	21
1.6	Evaluating the Results of Equating	22
1.7	Testing Situations Considered	23
1.8	Preview	24
1.9	Exercises	25
	References	26
2	Observed Score Equating Using the Random Groups Design	29
2.1	Mean Equating	30
2.2	Linear Equating	31

2.3	Properties of Mean and Linear Equating	32
2.4	Comparison of Mean and Linear Equating.	33
2.5	Equipercentile Equating.	36
2.5.1	Graphical Procedures	38
2.5.2	Analytic Procedures	42
2.5.3	Properties of Equated Scores in Equipercentile Equating	45
2.6	Estimating Observed Score Equating Relationships.	46
2.7	Scale Scores	50
2.7.1	Linear Conversions	50
2.7.2	Truncation of Linear Conversions.	53
2.7.3	Nonlinear Conversions	54
2.8	Equating Using Single Group Designs	60
2.9	Equating Using Alternate Scoring Schemes	60
2.10	Preview of What Follows	61
2.11	Exercises	62
	References	63
3	Random Groups: Smoothing in Equipercentile Equating	65
3.1	A Conceptual Statistical Framework for Smoothing	66
3.2	Properties of Smoothing Methods.	69
3.3	Presmoothing Methods	70
3.3.1	Polynomial Log-Linear Method	70
3.3.2	Strong True Score Method.	72
3.3.3	Illustrative Example	74
3.4	Postsmoothing Methods.	80
3.4.1	Illustrative Example	85
3.5	The Kernel Method of Equating.	89
3.6	Practical Issues in Equipercentile Equating	93
3.6.1	Summary of Smoothing Strategies	94
3.6.2	Smoothing and Population Distribution Irregularities.	95
3.6.3	Equating Error, Sample Size, and Smoothing Method	96
3.7	Exercises	98
	References	99
4	Nonequivalent Groups: Linear Methods	103
4.1	Tucker Method.	105
4.1.1	Linear Regression Assumptions	105
4.1.2	Conditional Variance Assumptions	106
4.1.3	Intermediate Results	107
4.1.4	Final Results	108
4.1.5	Special Cases	109

4.2	Levine Observed Score Method	109
4.2.1	Correlational Assumptions	110
4.2.2	Linear Regression Assumptions	110
4.2.3	Error Variance Assumptions.	111
4.2.4	Intermediate Results	111
4.2.5	General Results	112
4.2.6	Classical Congeneric Model Results	113
4.3	Levine True Score Method	116
4.3.1	Results	117
4.3.2	First-Order Equity.	119
4.4	Chained Linear Equating	121
4.4.1	Chained Linear Observed Score Equating	122
4.4.2	Chained Linear True Score Equating.	123
4.5	Illustrative Example and Other Topics	124
4.5.1	Illustrative Example	125
4.5.2	Synthetic Population Weights.	128
4.5.3	Mean Equating	128
4.5.4	Decomposing Observed Differences in Means and Variances.	129
4.5.5	Relationships Among Linear Observed Score Methods	132
4.5.6	Relationships Involving Levine Methods.	135
4.5.7	Other Issues Involving Methods	137
4.5.8	Scale Scores.	137
4.6	Appendix: Proof that $\sigma_s^2(T_X) = \gamma_1^2 \sigma_s^2(T_Y)$ Under the Classical Congeneric Model.	139
4.7	Exercises	139
	References	141
5	Nonequivalent Groups: Equipercentile Methods	143
5.1	Frequency Estimation Method	143
5.1.1	Conditional Distributions	144
5.1.2	Assumptions and Procedures	144
5.1.3	Numerical Example.	147
5.1.4	Estimating the Distributions.	150
5.1.5	Special Case: Braun-Holland Linear Method	151
5.1.6	Illustrative Example	152
5.2	Other Methods	158
5.2.1	Modified Frequency Estimation	158
5.2.2	Chained Equipercentile Equating	159
5.2.3	Illustrative Example	164
5.3	Practical Issues.	165
5.4	Exercises	166
	References	166

6	Item Response Theory Methods	171
6.1	Some Necessary IRT Concepts.	172
6.1.1	Unidimensionality and Local Independence Assumptions.	172
6.1.2	IRT Models	173
6.1.3	IRT Parameter Estimation	176
6.2	Transformations of IRT Scales	177
6.2.1	Transformation Equations	177
6.2.2	Demonstrating the Appropriateness of Scale Transformations	178
6.2.3	Expressing <i>A</i> and <i>B</i> Constants	179
6.2.4	Expressing <i>A</i> and <i>B</i> Constants in Terms of Groups of Items and/or Persons	180
6.3	Transforming IRT Scales When Parameters are Estimated.	181
6.3.1	Designs	182
6.3.2	Mean/Sigma and Mean/Mean Transformation Methods	183
6.3.3	Characteristic Curve Transformation Methods	184
6.3.4	Comparisons Among Scale Transformation Methods	189
6.4	Equating and Scoring	191
6.5	Equating True Scores	192
6.5.1	Test Characteristic Curves	192
6.5.2	True Score Equating Process	193
6.5.3	The Newton-Raphson Method	193
6.5.4	Using True Score Equating with Observed Scores	196
6.6	Equating Observed Scores	197
6.7	IRT True Score Versus IRT Observed Score Equating	201
6.8	Illustrative Example	201
6.8.1	Item Parameter Estimation and Scaling	202
6.8.2	IRT True Score Equating.	206
6.8.3	IRT Observed Score Equating	207
6.8.4	Rasch Equating	213
6.9	Using IRT Calibrated Item Pools and Other Designs	214
6.9.1	Common-Item Equating to a Calibrated Pool	215
6.9.2	Item Preequating.	219
6.9.3	Other Designs	221
6.10	Equating with Polytomous IRT	221
6.10.1	Polytomous IRT Models for Ordered Responses.	222
6.10.2	Scoring Function, Item Response Function, and Test Characteristic Curve.	227
6.10.3	Parameter Estimation and Scale Transformation with Polytomous IRT Models.	228

6.10.4	True Score Equating	232
6.10.5	Observed Score Equating.	232
6.10.6	Example Using the Graded Response Model	233
6.11	Robustness to Violations of the Unidimensionality Assumption	235
6.12	Practical Issues and Caveat	238
6.13	Exercises	239
	References	241
7	Standard Errors of Equating	247
7.1	Definition of Standard Error of Equating.	248
7.2	The Bootstrap	250
7.2.1	Standard Errors Using the Bootstrap	250
7.2.2	Standard Errors of Equating.	252
7.2.3	Parametric Bootstrap	253
7.2.4	Standard Errors of Equipercentile Equating with Smoothing	255
7.2.5	Standard Errors of Scale Scores	256
7.2.6	Standard Errors of Equating Chains	257
7.2.7	Mean Standard Error of Equating	258
7.2.8	Caveat.	259
7.3	The Delta Method	259
7.3.1	Mean Equating Using Single Group and Random Groups Designs	260
7.3.2	Linear Equating Using the Random Groups Design	261
7.3.3	Equipercentile Equating Using the Random Groups Design	263
7.3.4	Standard Errors for Other Designs	264
7.3.5	Illustrative Example	265
7.3.6	Approximations	267
7.3.7	Standard Errors for Scale Scores	268
7.3.8	Standard Errors of Equating Chains	269
7.3.9	Using Delta Method Standard Errors.	270
7.4	Using Standard Errors in Practice.	276
7.5	Exercises	278
	References	279
8	Practical Issues in Equating	283
8.1	Equating and the Test Development Process	285
8.1.1	Test Specifications	285
8.1.2	Changes in Test Specifications	286
8.1.3	Characteristics of Common-Item Sets	287

- 8.2 Data Collection: Design and Implementation 289
 - 8.2.1 Choosing Among Equating Designs 289
 - 8.2.2 Developing Equating Linkage Plans 292
 - 8.2.3 Examinee Groups Used in Equating 300
 - 8.2.4 Sample Size Requirements 303
- 8.3 Choosing from Among the Statistical Procedures 305
- 8.4 Equating Criteria and Designs in Research Studies 310
 - 8.4.1 Criteria and Designs Based on Error in Estimating Equating Relationships 310
 - 8.4.2 Equating in a Circle 318
 - 8.4.3 Criteria and Designs Based on Assessing Group Invariance of Equating Relationships 319
 - 8.4.4 Criteria and Designs Based on the Equity Property of Equating 320
 - 8.4.5 Discussion of Equating Criteria and Designs 325
- 8.5 Choosing from Among Equating Results in Operational Equating 326
 - 8.5.1 Equating Versus Not Equating 326
 - 8.5.2 Use of Robustness Checks 327
 - 8.5.3 Choosing from Among Results in the Random Groups Design 327
 - 8.5.4 Choosing from Among Results in the Common-Item Nonequivalent Groups Design 328
 - 8.5.5 Use of Consistency Checks 329
 - 8.5.6 Equating and Score Scales 330
- 8.6 Importance of Standardization Conditions and Quality Control Procedures 331
 - 8.6.1 Test Development 331
 - 8.6.2 Test Administration and Standardization Conditions 331
 - 8.6.3 Quality Control 333
 - 8.6.4 Reequating 334
- 8.7 Conditions Conducive to Satisfactory Equating 337
- 8.8 Comparability Issues in Special Circumstances 337
 - 8.8.1 Comparability Issues with Computer-Based Tests 337
 - 8.8.2 Comparability for Constructed-Response and Mixed-Format Tests 344
 - 8.8.3 Score Comparability with Optional Test Sections 348
- 8.9 Conclusion 349
- 8.10 Exercises 350
- References 352

9	Score Scales	371
9.1	Scaling Perspectives	372
9.2	Unit Scores, Item Scores, and Raw Scores.	377
9.2.1	Test Score Terminology	377
9.2.2	Unit and Item Scores	378
9.2.3	Raw Scores (Y)	380
9.3	Scores on Mixed-Format Tests	387
9.3.1	Weights Based on Numbers of Score Points	388
9.3.2	Observed Score Effective Weights	389
9.3.3	True Score Effective Weights.	390
9.3.4	Weights Chosen to Maximize Reliability.	390
9.3.5	Weighting Example.	391
9.3.6	Some Other Weighting Criteria and Issues.	392
9.3.7	Weights in IRT	392
9.4	Score Transformations.	393
9.5	Incorporating Normative Information	394
9.5.1	Linear Transformations	394
9.5.2	Nonlinear Transformations.	395
9.5.3	Example: Normalized Scale Scores.	397
9.5.4	Importance of Norm Group in Setting the Score Scale.	401
9.6	Incorporating Score Precision Information.	401
9.6.1	Rules of Thumb for Number of Distinct Score Points.	402
9.6.2	Linearly Transformed Score Scales with a Given Standard Error of Measurement	404
9.6.3	Score Scales with Approximately Equal Conditional Standard Errors of Measurement	405
9.6.4	Example: Incorporating Score Precision	407
9.6.5	Evaluating Psychometric Properties of Scale Scores.	410
9.6.6	The IRT θ -Scale as a Score Scale.	413
9.7	Incorporating Content Information	414
9.7.1	Item Mapping.	414
9.7.2	Scale Anchoring	415
9.7.3	Standard Setting	417
9.7.4	Numerical Example.	418
9.7.5	Practical Usefulness	420
9.8	Maintaining Score Scales.	420
9.9	Scales for Test Batteries and Composites	422
9.9.1	Test Batteries	422
9.9.2	Composite Scores	423
9.9.3	Maintaining Scales for Batteries and Composites	424

9.10 Vertical Scaling and Developmental Score Scales 425

 9.10.1 Structure of Batteries 427

 9.10.2 Type of Domain Being Measured 428

 9.10.3 Definition of Growth 429

 9.10.4 Designs for Data Collection for Vertical Scaling 431

 9.10.5 Test Scoring 434

 9.10.6 Hieronymus Statistical Methods 435

 9.10.7 Thurstone Statistical Methods 437

 9.10.8 IRT Statistical Methods 440

 9.10.9 Thurstone Illustrative Example 445

 9.10.10 IRT Illustrative Example 454

 9.10.11 Statistics for Comparing Scaling Results 461

 9.10.12 Some Limitations of Vertically Scaled Tests 463

 9.10.13 Vertical Scaling Designs with Variable Sections 465

 9.10.14 Maintaining Vertical Scales 466

 9.10.15 Research on Vertical Scaling 466

 9.10.16 Score Scales and Growth Models 471

9.11 Exercises 473

References 475

10 Linking 487

 10.1 Linking Categorization Schemes and Criteria 488

 10.1.1 Types of Linking 491

 10.1.2 Mislevy/Linn Taxonomy 492

 10.1.3 Holland and Dorans Framework 496

 10.1.4 Degrees of Similarity 498

 10.1.5 Summary and Other Approaches 500

 10.2 Group Invariance 501

 10.2.1 Statistical Methods Using Observed Scores 501

 10.2.2 Statistics for Overall Group Invariance 505

 10.2.3 Statistics for Pairwise Group Invariance 507

 10.2.4 Example: ACT and ITED Science Tests 508

 10.3 Additional Examples 527

 10.3.1 Extended Time 528

 10.3.2 Test Adaptations and Translated Tests 529

 10.4 Discussion 531

 10.5 Exercises 532

 References 533

Appendix A: Answers to Exercises 537

Appendix B: Computer Programs 559

Index 561

Notation

Arabic

1	Population taking Form X (Chapter 4)
2	Population taking Form Y (Chapter 4)
A	Slope constant in linear equating and raw-to-scale score transformations (Chapter 4)
A	Slope constant in IRT θ scale transformation (Chapter 6)
a	Item slope parameter in IRT (Chapter 6)
B	Location constant in linear equating and raw-to-scale score transformations (Chapter 4)
B	Location constant in IRT θ scale transformation (Chapter 6)
b	Item location parameter in IRT (Chapter 6)
b	Item or category location parameter in polytomous IRT (Chapter 6)
b*	Nonlinear transformation of b (Chapter 9)
bias	Bias (Chapter 3)
C	Number of degrees of the polynomial in log-linear smoothing (Chapter 3)
c	Item pseudochance level parameter in IRT (Chapter 6)
c	Item location parameter in Bock's nominal categories model (Chapter 6)
constant	A constant (Chapter 2)
cov	Sampling covariance (Chapter 7)
D	Scaling constant in IRT, usually set to 1.7 (Chapter 6)
DTM	Difference That Matters (Chapter 10)
d	Category location parameter in generalized partial credit model (Chapter 6)
$d_Y(x)$	Expected value of a cubic spline estimator of $e_Y(x)$ (Chapter 3)
$d^*_Y(x)$	Average of two splines (Chapter 3)
df	Degrees of freedom (Chapter 3)
E	Expected value (Chapter 1)
E	Number correct error score (Chapter 4)
e	The equipercntile equating function, such as $e_Y(x)$ (Chapter 2)

$e_Y(x)$	The Form Y equipercentile equivalent of a Form X score (Chapter 1)
$e_X(y)$	The Form X equipercentile equivalent of a Form Y score (Chapter 2)
<i>effect size</i>	Effect size (Chapter 9)
<i>eq</i>	General equating function, such as $eq_Y(x)$ (Chapter 1)
<i>ew</i>	Effective weight (Chapter 9)
<i>ewMAD</i>	Equally weighted average of absolute differences (Chapter 10)
<i>ewMD</i>	Equally weighted average of differences (Chapter 10)
<i>ewREMSD</i>	Equally weighted Root Expected Mean Square Difference (Chapter 10)
<i>exp</i>	Exponential (Chapter 6)
$F(x)$	$Pr(X \leq x)$ is the cumulative distribution for X (Chapter 1)
F^*	Cumulative distribution function of $eq_X(y)$ (Chapter 2)
F^{-1}	Inverse of function F (Chapter 2)
f	A general function (Chapter 7)
f'	The first derivative of f (Chapter 7)
$f(x)$	$Pr(X = x)$ is the discrete density for X (Chapter 2)
$f(x, v)$	$Pr(X = x \text{ and } V = v)$ is the joint density of X and V (Chapter 5)
$f(x v)$	$Pr(X = x \text{ given } V = v)$ is the conditional density of x given v (Chapter 5)
<i>func</i>	Function solved for in Newton–Raphson iterations (Chapter 6)
<i>func'</i>	First derivative of function solved for in Newton–Raphson iterations (Chapter 6)
$G(y)$	$Pr(Y \leq y)$ is the cumulative distribution for Y (Chapter 1)
G^*	The cumulative distribution function of e_Y (Chapter 1)
G^{-1}	Inverse of function G (Chapter 2)
g	Item subscript in IRT (Chapter 6)
g	Index used to sum over categories in generalized partial credit model (Chapter 6)
g	Arcsine transformed proportion-correct score (Chapter 9)
$g(y)$	$Pr(Y = y)$ is the discrete density for Y (Chapter 2)
$g(y, v)$	$Pr(Y = y \text{ and } V = v)$ is the joint density of Y and V (Chapter 5)
$g(y v)$	$Pr(Y = y \text{ given } V = v)$ is the conditional density of y given v (Chapter 5)
g_{adj}	Density adjusted by adding 10^{-6} to each density and then standardizing (Chapter 2)
H	Number of subgroups (Chapter 10)
H_{crit}	Criterion function for Haebara's method (Chapter 6)
H_{diff}	Difference function for Haebara's method (Chapter 6)
h	Index for summing over categories (Chapter 6)
h	Number of scale score points for a confidence interval (Chapter 9)
h	Subgroup designator (Chapter 10)
$h(v)$	$Pr(V = v)$ is the discrete density for V (Chapter 5)
I	IRT scale (Chapter 6)
I	Number of scale scores on Test X (Chapter 10)
i and i'	Individuals (Chapter 6)

<i>intercept</i>	Intercept of an equating function (Chapter 2)
<i>irt</i>	IRT true-score equating function (Chapter 6)
<i>J</i>	IRT scale (Chapter 6)
<i>J</i>	Number of scale scores on Test Y (Chapter 10)
<i>j</i> and <i>j'</i>	Items (Chapter 6)
<i>K</i>	Number of items (Chapter 2)
<i>KR-20</i>	Kuder–Richardson Formula 20 reliability coefficient (Chapter 9)
<i>KR-21</i>	Kuder–Richardson Formula 21 reliability coefficient (Chapter 9)
<i>k</i>	Lord's <i>k</i> in the Beta4 method (Chapter 3)
<i>k</i>	Categories for an item in polytomous IRT (Chapter 6)
<i>ku</i>	Kurtosis, such as $ku(X) = E[X - \mu(X)]^4 \sigma^4(X)$ (Chapter 2)
$I_Y(x)$	The Form Y linear equivalent of a Form X score (Chapter 2)
$I_X(y)$	The Form X linear equivalent of a Form Y score (Chapter 2)
<i>MAD</i>	Weighted average of absolute differences (Chapter 10)
<i>MD</i>	Weighted average of differences (Chapter 10)
<i>m</i>	Number of categories for an item in polytomous IRT (Chapter 6)
$m_Y(x)$	The mean equating equivalent of a Form X score (Chapter 2)
$m_X(y)$	The mean equating equivalent of a Form Y score (Chapter 2)
<i>max</i>	Maximum score (Chapter 6)
<i>min</i>	Minimum score (Chapter 6)
<i>mse</i>	Mean squared error (Chapter 3)
<i>N</i>	Number of examinees (Chapter 2)
<i>NCE</i>	Normal Curve Equivalent (unrounded) (Chapter 9)
<i>NCE_{int}</i>	Normal Curve Equivalent rounded to an integer (Chapter 9)
$P(x)$	The percentile rank function for X (Chapter 2)
<i>P*</i>	A given percentile rank (Chapter 2)
<i>P**</i>	$P/100$ (Chapter 7)
P^{-1}	The percentile function for X (Chapter 2)
<i>p</i>	Probability of a correct response in IRT (Chapter 6)
<i>p</i>	Category response function in polytomous IRT (Chapter 6)
<i>p*</i>	Cumulative category response function in polytomous IRT (Chapter 6)
<i>p'</i>	First derivative of <i>p</i> (Chapter 6)
pl_{Yh}	Parallel linear equating equivalent on Test Y for subgroup <i>h</i> (Chapter 10)
$Q(y)$	Percentile rank function for Y (Chapter 2)
Q^{-1}	Percentile function for Y (Chapter 2)
<i>R</i>	Number of bootstrap replications (Chapter 7)
<i>REMSD</i>	Root Expected Mean Square Difference (Chapter 10)
<i>RMSD</i>	Root Mean Square Difference (Chapter 10)
<i>RP</i>	Response Probability level in item mapping (Chapter 9)
<i>r</i>	Index for calculating observed score distribution in IRT (Chapter 6)
<i>r</i>	Index for bootstrap replications (Chapter 7)
<i>rmsel</i>	Root mean squared error for linking (Chapter 10)

S	Smoothing parameter in postsMOOTHING (Chapter 3)
SC	Scale score random variable (Chapter 9)
$SLcrit$	Criterion function for Stocking-Lord method (Chapter 6)
$SLdiff$	Difference function for Stocking-Lord method (Chapter 6)
SMD	Standardized Mean Difference (Chapter 10)
s	Synthetic population (Chapter 4)
sc	Scale score transformation, such as $sc(y)$ (Chapter 2)
sc_{int}	Scale score rounded to an integer (Chapter 2)
se	Standard error, such as $se(x)$ (Chapter 3)
sem	Standard error of measurement (Chapter 7)
sk	Skewness, such as $sk(X) = E[X - \mu(X)]^3 / \sigma^3(X)$ (Chapter 2)
$slope$	Slope of equating function (Chapter 2)
st	Stanine (unrounded) (Chapter 9)
st	Scaling test (Chapter 9)
st_{int}	Stanine rounded to an integer (Chapter 9)
T	Number correct true score (Chapter 4)
T	Normalized score with mean of 50 and standard deviation of 10 (Chapter 9)
T_{int}	Normalized score with mean of 50 and standard deviation of 10 rounded to an integer (Chapter 9)
t	Realization of number correct true score (Chapter 4)
$t_Y(x)$	Expected value of an alternate estimator of $e_Y(x)$ (Chapter 3)
U	Uniform random variable (Chapter 2)
u	Standard deviation units (Chapter 7)
V	The random variable indicating raw score on Form V (Chapter 4)
v	Spline coefficient (Chapter 3)
v	A realization of V (Chapter 4)
v	Subgroup weight for a particular score (Chapter 10)
var	Sampling variance (Chapter 3)
w	Weight for synthetic group (Chapter 4)
w	Nominal weight (Chapter 9)
w	Subgroup weight (Chapter 10)
X	The random variable indicating raw score on Form X (Chapter 1)
X	Random variable indicating scale score on Test X (Chapter 10)
X^*	Equals $X + U$, used in the continuization process (Chapter 2)
x	A realization of X (Chapter 2)
x^*	Integer closest to x such that $x^* - .5 \leq x < x^* + .5$ (Chapter 2)
x^*	Form X_2 score equated to the Form X_1 scale (Chapter 7)
x_{high}	Upper limit in spline calculations (Chapter 3)
x_L^*	The largest integer score with a cumulative percent less than P^* (Chapter 2)
x_{low}	Lower limit in spline calculations (Chapter 3)
x_U^*	Smallest integer score with a cumulative percent greater than P^* (Chapter 2)

Y	The random variable indicating raw score on Form Y (Chapter 1)
Y	Random variable indicating scale score on Test Y (Chapter 10)
y	A realization of Y (Chapter 1)
y_i^*	Largest tabled raw score less than or equal to $e_Y(x)$ in finding scale scores (Chapter 2)
y_L^*	The largest integer score with a cumulative percent less than Q^* (Chapter 2)
y_U^*	The smallest integer score with a cumulative percent greater than Q^* (Chapter 2)
Z	The random variable indicating raw score on Form Z (Chapter 4)
z	A realization of Z (Chapter 4)
z	Unit normal variable (Chapter 7)
z	Normalized score (Chapter 10)
z^*	Selected set of normalized scores in Thurstone scaling (Chapter 9)
z_γ	Unit normal score associated with a 100γ % confidence interval (Chapter 9)

Greek

$\alpha(X V)$	Linear regression slope (Chapter 4)
$\alpha(Y V)$	Linear regression slope (Chapter 4)
$\beta(X V)$	Linear regression intercept (Chapter 4)
$\beta(Y V)$	Linear regression intercept (Chapter 4)
χ^2	Chi-square test statistic (Chapter 3)
δ	Location parameter in congeneric models (Chapter 4)
ϕ	Normal ordinate (Chapter 7)
γ	Expansion factor in linear equating with the common-item nonequivalent groups design (Chapter 4)
γ	Confidence coefficient (Chapter 9)
λ	Effective test length in congeneric models (Chapter 4)
μ	Mean as in $\mu(X)$ and $\mu(Y)$ (Chapter 2)
ν	Weight for a pair of subgroups and a particular score (Chapter 10)
Φ	Inverse normal transformation (Chapter 9)
Θ	Parameter used in developing the delta method (Chapter 7)
θ	Ability in IRT (Chapter 6)
θ^+	New value in Newton–Raphson iterations (Chapter 6)
θ_-	Initial value in Newton–Raphson iterations (Chapter 6)
θ^*	Nonlinear transformation of θ (Chapter 9)
ρ	Correlation, such as $\rho(X, V)$ (Chapter 4)
$\rho(X, X')$	Reliability (Chapter 4)
$\sigma(X, V)$	Covariance between X and V (Chapter 4)
$\sigma(Y, V)$	Covariance between Y and V (Chapter 4)
σ^2	Variance such as $\sigma^2(X) = \mathbf{E}[X - \mu(X)]^2$ (Chapter 4)
σ_{ij}	Covariance between variables i and j (Chapter 9)

τ	True score (Chapter 1)
τ^*	True score outside of range of possible true scores (Chapter 6)
$\hat{\tau}$	Estimated true scores (Chapter 9)
ω	Weight in log-linear smoothing (Chapter 3)
Ψ	Function that relates true scores (Chapter 4)
ψ	Distribution of a latent variable (Chapter 3)
∂	Partial derivative (Chapter 7)

Chapter 1

Introduction and Concepts

This chapter provides a general overview of equating and briefly considers important concepts. The concept of equating is described, as is why it is needed, and how to distinguish it from other related processes. Equating properties and designs are considered in detail, because these concepts provide the organizing themes for addressing the statistical methods treated in subsequent chapters. Some issues in evaluating equating are also considered. The chapter concludes with a preview of subsequent chapters.

1.1 Equating and Related Concepts

Scores on tests often are used as one piece of information in making important decisions. Some of these decisions focus at the *individual level*, such as when a student decides which college to attend or on a course in which to enroll. For other decisions the focus is more at an *institutional level*. For example, an agency or institution might need to decide what test score is required to certify individuals for a profession or to admit students into a college, university, or the military. Still other decisions are made at the *public policy level*, such as addressing what can be done to improve education in the United States and how changes in educational practice can be evaluated. Regardless of the type of decision that is to be made, it should be based on the most accurate information possible: All other things being equal, *the more accurate the information, the better the decision*.

Making decisions in many of these contexts requires that tests be administered on multiple occasions. For example, college admissions tests typically are administered on particular days, referred to as *test dates*, so examinees can have some flexibility in choosing when to be tested. Tests also are given over many years to track educational trends over time. If the same test questions were routinely administered on each test

Some of the material in this chapter is based on Kolen (1988).

date, then examinees might inform others about the test questions. Or, an examinee who tested twice might be administered the same test questions on the two test dates. In these situations, a test might become more of a measure of exposure to the specific questions that are on the test than of the construct that the test is supposed to measure.

1.1.1 Test Forms and Test Specifications

These test security problems can be addressed by administering a different collection of test questions, referred to as a *test form*, to examinees who test on different test dates. A test form is a set of test questions that is built according to content and statistical *test specifications* (Schmeiser and Welch 2006). Test specifications provide guidelines for developing the test. Those responsible for constructing the test, the *test developers*, use these specifications to ensure that the test forms are as similar as possible to one another in content and statistical characteristics.

1.1.2 Equating

The use of different test forms on different test dates leads to another concern: the forms might differ somewhat in difficulty. *Equating* is a statistical process that is used to adjust scores on test forms so that scores on the forms can be used interchangeably. Equating adjusts for differences in difficulty among forms that are built to be similar in difficulty and content.

The following situation is intended to develop further the concept of equating. Suppose that a student takes a college admissions test for the second time and earns a higher reported score than on the first testing. One explanation of this difference is that the reported score on the second testing reflects a higher level of achievement than the reported score on the first testing. However, suppose that the student had been administered exactly the same test questions on both testings. Rather than indicating a higher level of achievement, the student's reported score on the second testing might be inflated because the student had already been exposed to the test items. Fortunately, a new test form is used each time a test is administered for most college admissions tests. Therefore, a student would not likely be administered the same test questions on any two test dates.

The use of different test forms on different test dates might cause another problem, as is illustrated by the following situation. Two students apply for the same college scholarship that is based partly on test scores. The two students take the test on different test dates, and Student 1 earns a higher reported score than Student 2. One possible explanation of this difference is that Student 1 is higher achieving than Student 2. However, if Student 1 took an easier test form than Student 2, then Student 1 would have an unfair advantage over Student 2. In this case, the difference in scores might be due to differences in the difficulty of the test forms rather than in

the achievement levels of the students. To avoid this problem, equating is used with most college admissions tests. If the test forms are successfully equated, then the difference in equated scores for Student 1 and Student 2 is not attributable to Student 1's taking an easier form.

The process of equating is used in situations where such *alternate forms* of a test exist and scores earned on different forms are compared to each other. Even though test developers attempt to construct test forms that are as similar as possible to one another in content and statistical specifications, the forms typically differ somewhat in difficulty. Equating is intended to adjust for these difficulty differences, allowing the forms to be used interchangeably. *Equating adjusts for differences in difficulty, not for differences in content.* After successful equating, for example, examinees who earn an equated score of, say, 26 on one test form could be considered, on average, to be at the same achievement level as examinees who earn an equated score of 26 on a different test form.

1.1.3 Processes That are Related to Equating

There are processes that are similar to equating, which may be more properly referred to as *scaling to achieve comparability*, in the terminology of the *Standards for Educational and Psychological Testing* (AERA, APA, NCME 1999), or *linking*, in the terminology of Holland and Dorans (2006), Linn (1993) and Mislevy (1992). One of these processes is *vertical scaling* (frequently referred to as *vertical "equating"*), which often is used with elementary school achievement test batteries. In these batteries, students often are administered questions that test content matched to their current grade level. This procedure allows developmental scores (e.g., grade equivalents) of examinees at different grade levels to be compared. Because the content of the tests administered to students at various educational levels is different, however, scores on tests intended for different educational levels cannot be used interchangeably. Other examples of linking include relating scores on one test to those on another, and scaling the tests within a battery so that they all have the same distributional characteristics. As with vertical scaling, solutions to these problems do not allow test scores to be used interchangeably, because the content of the tests is different.

Although similar statistical procedures often are used in linking and equating, their purposes are different. Whereas tests that are purposefully built to be different are linked, equating is used to adjust scores on test forms that are built to be as similar as possible in content and statistical characteristics. When equating is successful, scores on alternate forms can be used interchangeably. Issues in linking tests that are not built to the same specifications are considered further in Chaps. 9 and 10.

1.1.4 Equating and Score Scales

On a multiple-choice test, the *raw score* an examinee earns is often the number of items the examinee answers correctly. Other raw scores might involve penalties for wrong answers or weighting items differentially. On tests that require ratings by judges, a raw score might be the sum of the numerical ratings made by the judges.

Raw scores often are transformed to *scale scores*. The *raw-to-scale score transformation* can be chosen by test developers to enhance the interpretability of scores by incorporating useful information into the score scale (Kolen 2006; Petersen et al. 1989). Information based on a nationally representative group of examinees, referred to as a *national norm group*, sometimes is used as a basis for establishing score scales. For example, the number-correct scores for the four tests of the initial form of a revised version of the ACT tests were scaled (Brennan 1989) to have a mean scale score of 18 for a nationally representative sample of college-bound 12th graders. Thus, an examinee who earned a scale score of 22, for example, would know that this score was above the mean scale score for the nationally representative sample of college-bound 12th graders used to develop the score scale. One alternative to using nationally representative norm groups is to base scale score characteristics on a *user norm group*, which is a group of examinees that is administered the test under operational conditions. For example, a rescaled SAT scale was established for use beginning in 1995 by setting the mean score equal to 500 for the group of SAT examinees that graduated from high school in 1990 (Cook 1994; Dorans 2002).

Scaling and Equating Process

Equating can be viewed as an aspect of a more general *scaling and equating process*. Score scales typically are established using a single test form. For subsequent test forms, the scale is maintained through an equating process that places raw scores from subsequent forms on the established score scale. In this way, a scale score has the same meaning regardless of the test form administered or the group of examinees tested. Typically, raw scores on the new form are equated to raw scores on the old form, and these equated raw scores are then converted to scale scores using the raw-to-scale score transformation for the old form.

Example of the Scaling and Equating Process

The hypothetical conversions shown in Table 1.1 illustrate the scaling and equating process. The first two columns show the relationship between Form Y raw scores and scale scores. For example, a raw score of 28 on Form Y converts to a scale score of 14 (At this point there is no need to be concerned about what particular method was used to develop the raw-to-scale score transformation). The relationship between Form

Table 1.1 Hypothetical conversion tables for test forms

Scale	Form Y raw	Form X ₁ raw	Form X ₂ raw
•	•	•	•
•	•	•	•
•	•	•	•
13	26	27	28
14	27	28	29
14	28	29	30
15	29	30	31
15	30	31	32
•	•	•	•
•	•	•	•
•	•	•	•

Y raw scores and scale scores shown in the first two columns involves scaling—not equating, because Form Y is the only form that is being considered so far.

Assume that an equating process indicates that Form X₁ is 1 point easier than Form Y throughout the score scale. A raw score of 29 on Form X₁ would thus reflect the same level of achievement as a raw score of 28 on Form Y. This relationship between Form Y raw scores and Form X₁ raw scores is displayed in the second and third columns in Table 1.1. What scale score corresponds to a Form X₁ raw score of 29? A scale score of 14 corresponds to this raw score, because a Form X₁ raw score of 29 corresponds to a Form Y raw score of 28, and a Form Y raw score of 28 corresponds to a scale score of 14.

To carry the example one step further, assume that Form X₂ is found to be uniformly 1 raw score point easier than Form X₁. Then, as illustrated in Table 1.1, a raw score of 30 on Form X₂ corresponds to a raw score of 29 on Form X₁, which corresponds to a raw score of 28 on Form Y, which corresponds to a scale score of 14. Later, additional forms could be converted to scale scores by a similar chaining process. The result of a successful scaling and equating process is that scale scores on all forms can be used interchangeably.

Possible Alternatives to Equating

Equating has the potential to improve score reporting and interpretation of tests that have alternate forms when examinees administered different forms are evaluated at the same time, or when score trends are to be evaluated over time. When at least one of these characteristics is present, at least two possible, but typically unacceptable, alternatives to equating exist. One alternative is to report raw scores regardless of the form administered. As was the case with Students 1 and 2 considered earlier, this approach could cause problems because examinees who were administered an easier form are advantaged and those who were administered a more difficult